

Multivariabilna logistična regresija s ponovitvijo linearne regresije

doc. dr. Mitja Kos, mag. farm.

Katedra za socialno farmacijo
Univerza v Ljubljani- Fakulteta za farmacijo

Analiza povezanosti

□ Regresija:

- Statistična analiza povezanosti dveh spremenljivk pri čemer je ena odvisna druga pa neodvisna spremenljivka (njene vrednosti lahko sami izbiramo).
 - Opazovani pojav = odvisna spremenljivka
 - Napovedni dejavnik = neodvisna spremenljivka
 - Npr. Učinek = $f(\text{koncentracije})$, regresijski koeficienti b

□ Korelacija:

- Statistična analiza povezanosti dveh spremenljivk, ki sta obe naključni (njihove vrednosti ne moremo vnaprej izbirati) in odvisni (na obe delujejo biološki in drugi dejavniki variabilnosti)
 - Skušamo kvantificirati stopnjo povezanosti dveh spremenljivk. Npr. Pearsonov korelacijski koeficient R , determinacijski koeficient R^2

- Pri praktičnem delu ni ostre meje med obema metodama.
-

Analiza povezanosti

Statistični modeli:

■ Univariabilni:

en napovedni dejavnik

Povezava kot pomembna pokaže:

■ zaradi dejanske povezanosti napovednega dejavnika s pojavom

■ lahko tudi zaradi povezanosti z nekim drugim napovednim dejavnikom.

■ Multivariabilni:

več napovednih dejavnikov

Linearna regresija

- Preprosta LR: matematični model = premica
 - za vsak posamezni y

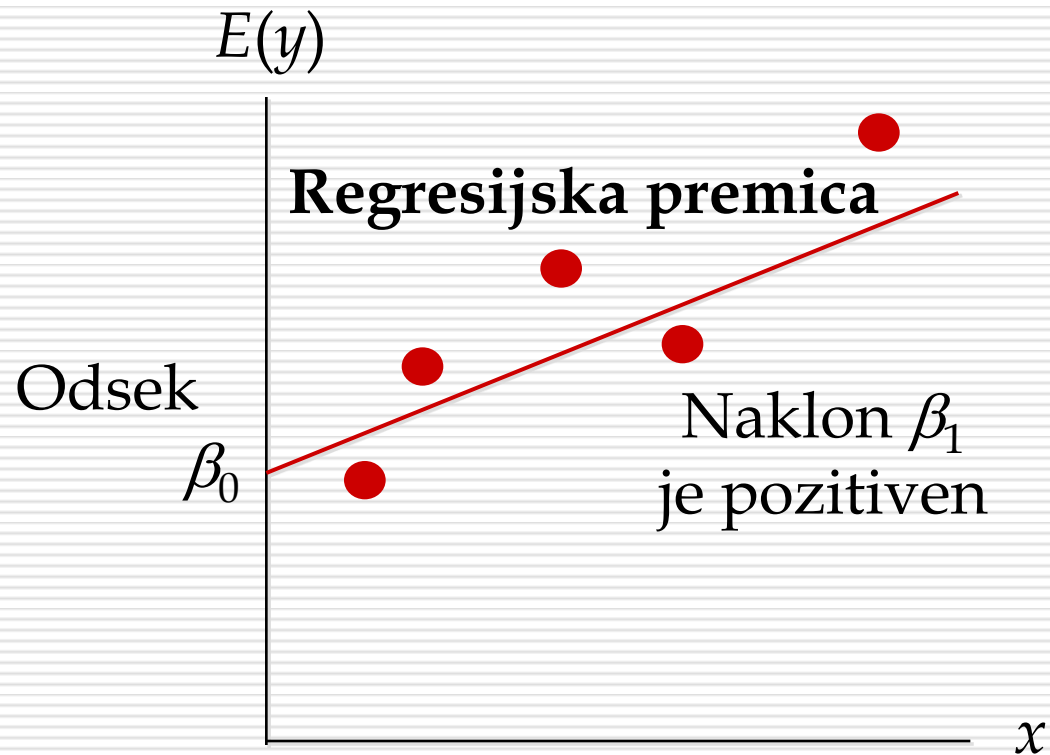
$$y = \beta_0 + \beta_1 x + \varepsilon$$

- β_0 in β_1 sta parametera modela.
 - ε je napaka $N(0, \sigma_e^2)$
- Pričakovana vrednost (povprečen y)

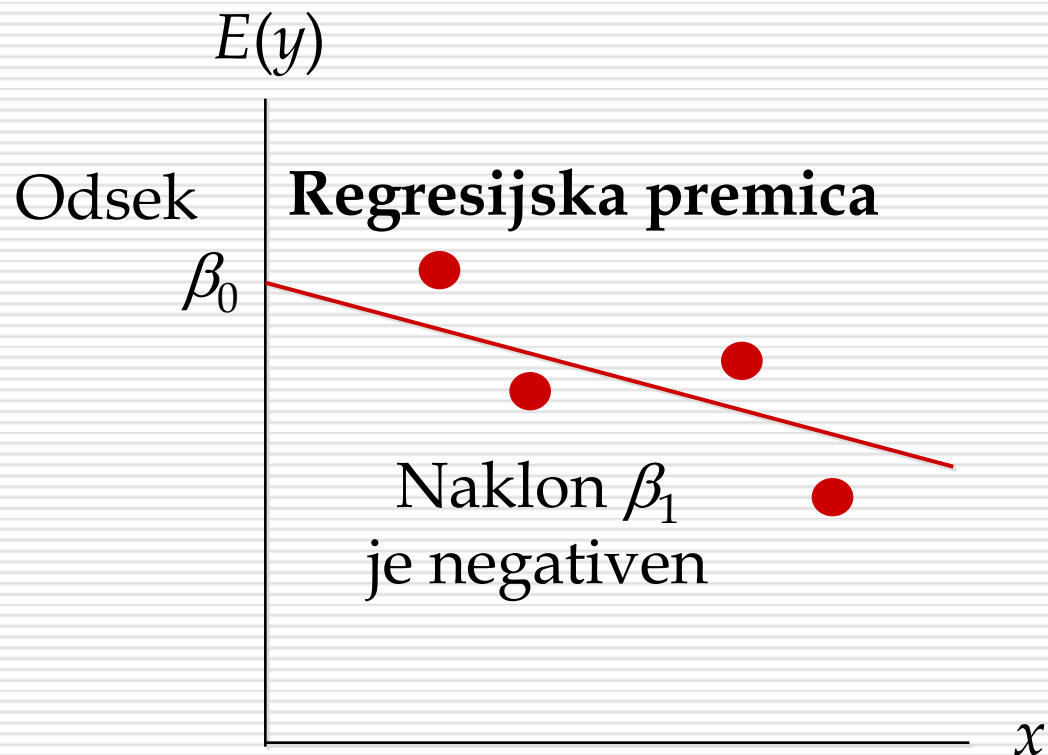
$$E(y) = \beta_0 + \beta_1 x$$

$$\beta_1 = \frac{\Delta y}{\Delta x}$$

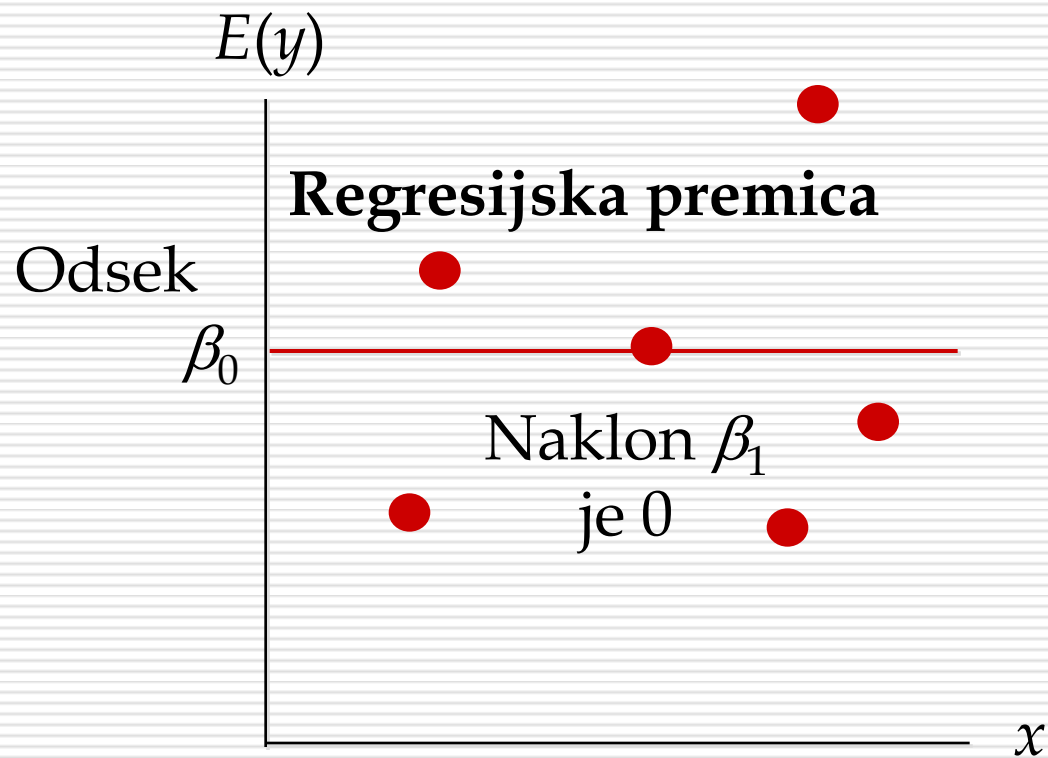
Pozitiven linearni odnos



Negativen linearni odnos



Ni povezave

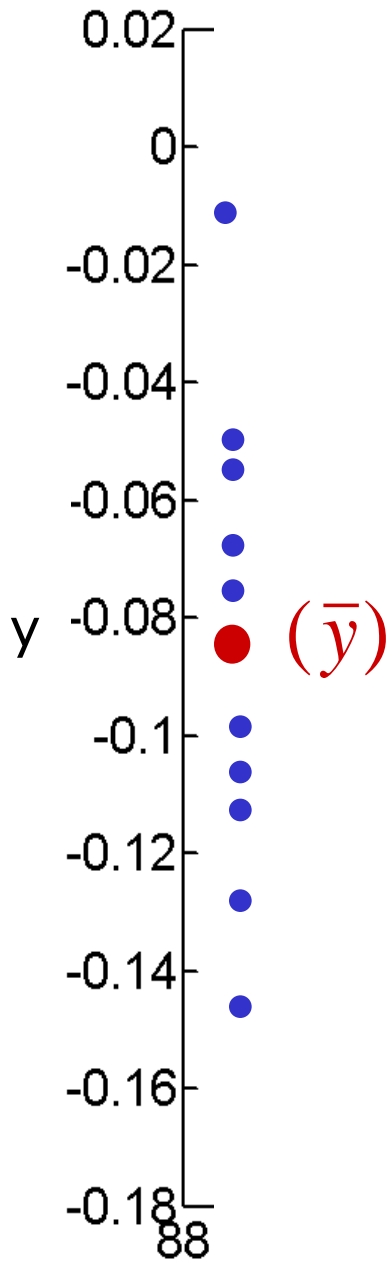


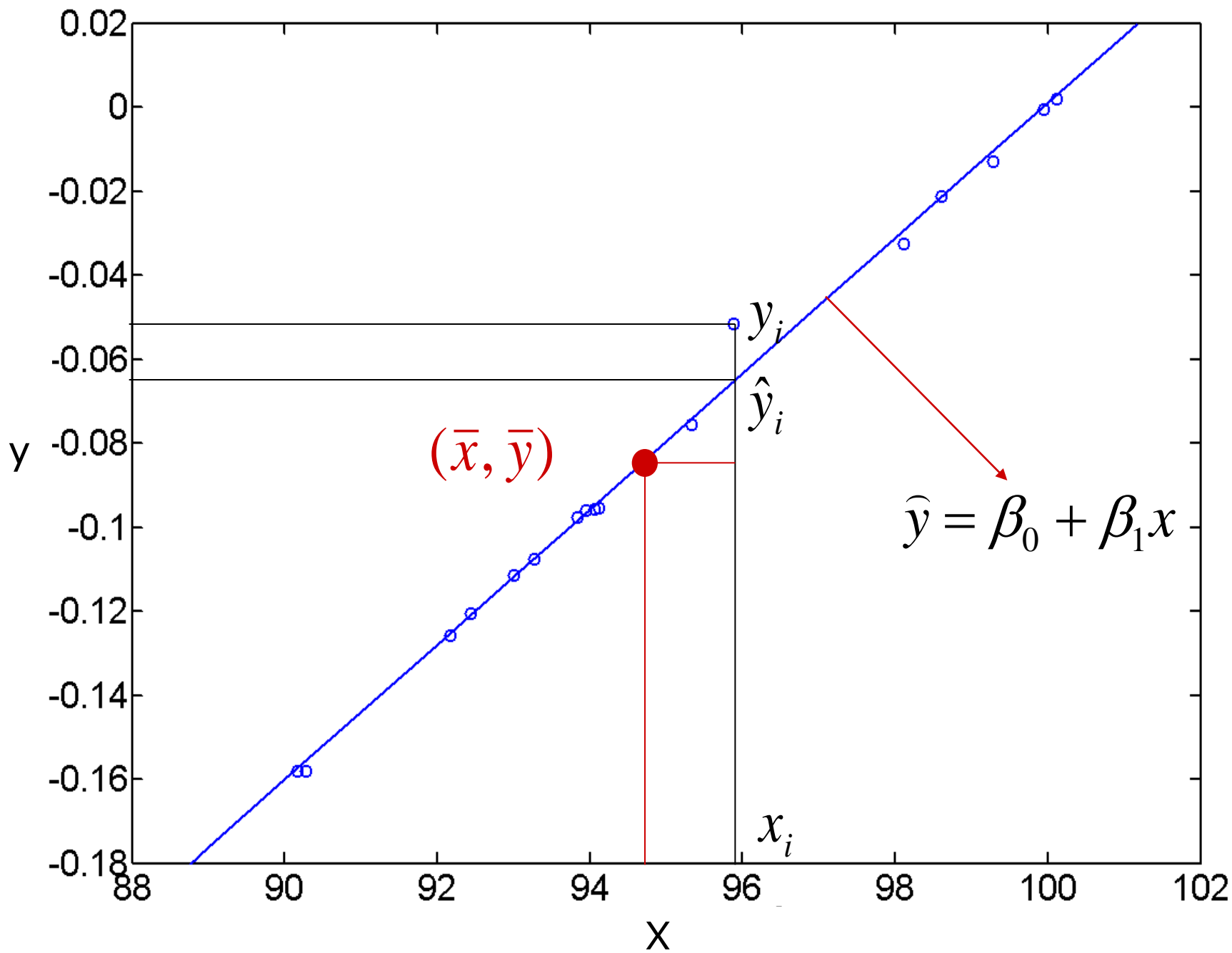
Metoda najmanjših kvadratov

$$\min \sum (y_i - \hat{y}_i)^2$$

\hat{y}_i = ocena i-te vrednosti odvisne spremenljivke

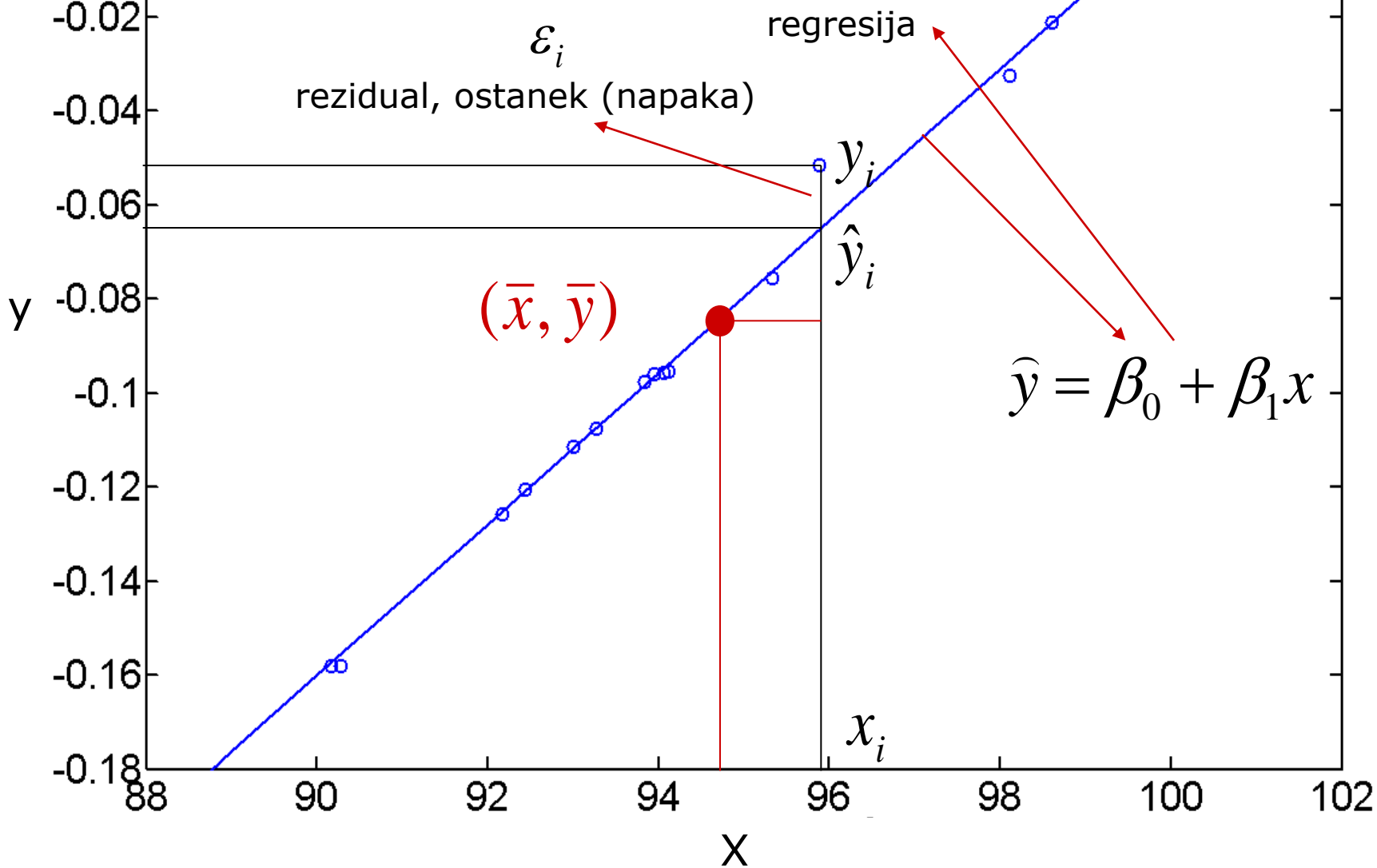
y_i = opažena i-ta vrednosti odvisne spremenljivke





$$\sum (y_i - \bar{y})^2 = \sum (\hat{y}_i - \bar{y})^2 + \sum (y_i - \hat{y}_i)^2$$

$$(y_i - \bar{y}) = (y_i - \hat{y}_i) - (\hat{y}_i - \bar{y})$$



Multipli regresijski model

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \varepsilon$$

- $\beta_0, \beta_1, \beta_2, \dots, \beta_p$ so parametri.
 - Standardizirani parametri ($x_i \rightarrow z_i$)
 - ε je napaka, ki je slučajna spremenljivka.
-

Metoda najmanjših kvadratov

□ Kriterij

$$\min \sum (y_i - \hat{y}_i)^2$$

□ Določitev koeficientov

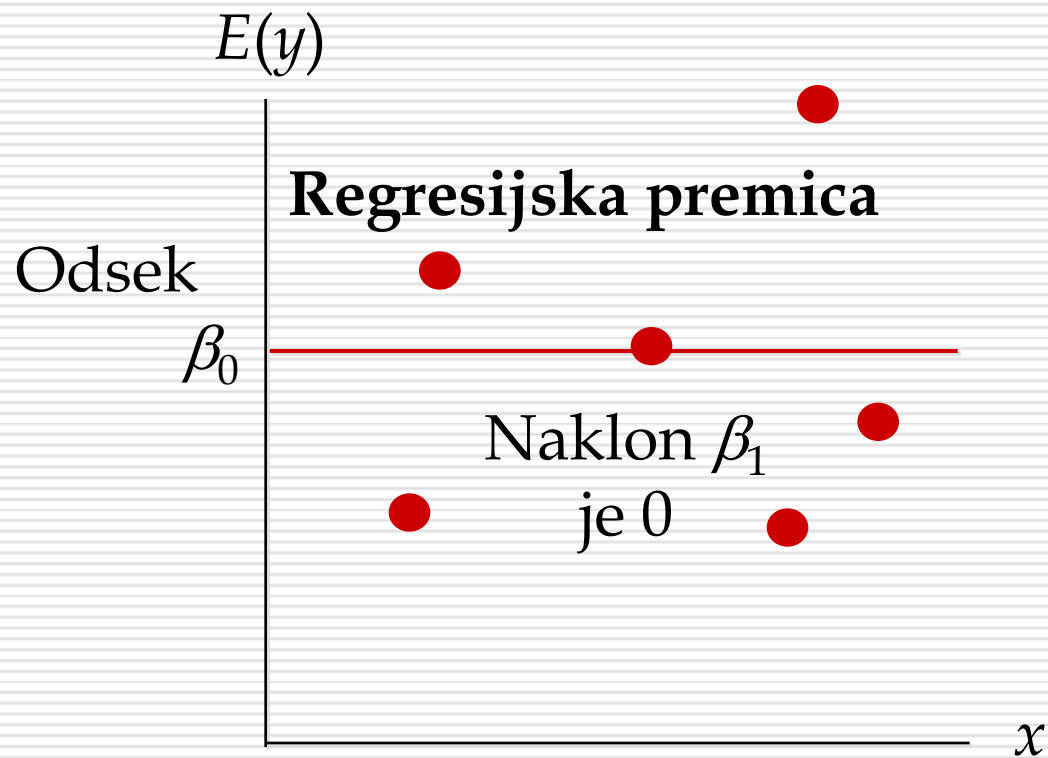
Kompleksna algebra. Uporaba statističnih programskih paketov.

□ Interpretacija koeficientov

$$E(y) = a + b_1x_1 + b_2x_2 + \dots + b_kx_k$$

$$b_k = \frac{\Delta y}{\Delta x_k}$$

Ni povezave



Statistično sklepanje

- Pri enostavni linearni regresiji nas F and t test vodita k istim sklepom.
 - V multipli regresiji uporabimo F in t test za različna namena.
-

F test

- Pomembnost modela kot celote.
 - Test for overall significance.
 - Hipotezi:
 - $H_0: \beta_1 = \beta_2 = \dots = \beta_p = 0$
 - H_a : Najmanj eden izmed parametrov ni enak 0
-

t test

- Če je izid F testa za model kot celota značilen, uporabimo t teste za ugotavljanje značilnosti vplivov posameznih neodvisnih spremenljivk.
- Za vsako neodvisno spremenljivko izvedemo en t test.
- Vsak t test imenujemo tudi test posamične značilnosti (test for individual significance).
- Hipotezi:

$$H_0: \beta_i = 0$$

$$H_a: \beta_i \neq 0$$

Koeficient determinacije

- Multipli koeficient determinacije R^2
 - Kako dobro model opisuje podatke?
-

Primer:

Vpliv kajenja na pljučno funkcijo

Z raziskavo želimo opredeliti vpliv kajenja na pljučno funkcijo. Naključno smo izbrali 1000 ljudi obeh spolov in spremljali pljučno funkcijo (parameter FEV), poleg tega smo beležili še njihovo starost in telesno višino. Za analizo rezultatov raziskave smo uporabili metodo multiple regresije. Neodvisne spremenljivke v analizi so bile:

starost [leta]

telesna višina [cm]

spol (0 = ženski, 1 = moški)

kajenje (0 = nekadilec, 1 = kadilec)

Odvisna spremenljivka pa je bila FEV [liter].

Primer:

Vpliv kajenja na pljučno funkcijo

Model Summary

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	.497 ^a	.247	.243	.25269

a. Predictors: (Constant), kajenje, t. višina (cm), starost (leta), spol

ANOVA^b

Model		Sum of Squares	df	Mean Square	F	Sig.
1	Regression	20.786	4	5.196	81.383	.000 ^a
	Residual	63.533	995	.064		
	Total	84.319	999			

a. Predictors: (Constant), kajenje, t. višina (cm), starost (leta), spol

b. Dependent Variable: FEV (l)

Coefficients^a

Model		Unstandardized Coefficients		Standardized Coefficients	t	Sig.
		B	Std. Error	Beta		
1	(Constant)	2.666	.118		22.568	.000
	starost (leta)	-1.9E-03	.001	-.063	-2.296	.022
	t. višina (cm)	-5.6E-04	.001	-.024	-.844	.399
	spol	.219	.017	.376	13.208	.000
	kajenje	-.275	.023	-.327	-11.857	.000

a. Dependent Variable: FEV (l)





TITANIC

www.titanicmovie.com

Tveganje, obeti ter razmerja

Spol	Izid		Skupaj
	Smrt	Preživetje	
Moški	$n_{11}=1364$	$n_{12}=367$	$n_{1+}=1731$
Ženski	$n_{21}=126$	$n_{22}=344$	$n_{2+}=470$
Skupaj	$n_{+1}=1490$	$n_{+2}=711$	$n=2201$

Tveganje, obeti ter razmerja

Relativno tveganje
ang. Relative Risk

$$RR = \frac{p_1}{p_2}$$

Razmerje obetov
ang. Odds Ratio

$$\theta = \frac{\frac{p_1}{1-p_1}}{\frac{p_2}{1-p_2}}$$

Povezava med razmerjem obetov ter relativnim tveganjem:

$$\theta = \frac{\frac{p_1}{1-p_1}}{\frac{p_2}{1-p_2}} = \frac{p_1}{p_2} \times \frac{1-p_2}{1-p_1} = RR \times \frac{1-p_2}{1-p_1}$$

Tveganje, obeti ter razmerja

$$p_1 = \frac{n_{11}}{n_{1+}} = \frac{1364}{1731} = 0,79$$

Spol	Izid		Skupaj
	Smrt	Preživetje	
Moški	$n_{11}=1364$	$n_{12}=367$	$n_{1+}=1731$
Ženski	$n_{21}=126$	$n_{22}=344$	$n_{2+}=470$
Skupaj	$n_{+1}=1490$	$n_{+2}=711$	$n=2201$

$$\theta = \frac{\frac{p_1}{1-p_1}}{\frac{p_2}{1-p_2}} = \frac{\frac{n_{11}/n_{1+}}{n_{12}/n_{1+}}}{\frac{n_{21}/n_{2+}}{n_{22}/n_{2+}}} = \frac{n_{11}n_{22}}{n_{12}n_{21}}$$

$$\theta_m = \frac{\frac{n_{11}}{n_{21}}}{\frac{n_{12}}{n_{22}}} = \frac{n_{11}n_{22}}{n_{12}n_{21}}$$

Preprosta logistična regresija

- Opazovani pojav = binarna spremenljivka 0/1
 - Kako uporabiti linearni model za opisovanje odnosa?
-

Preprosta logistična regresija

Vrednost binarne spremenljivke zapisati kot:

□ Verjetnost, da zavzame vrednost 1 pri danem X :

□ $\pi(x)$ - populacija

□ $p(x)$ - vzorec

■ Zavzame le vrednosti med 0 in 1- preozek razpon

□ Razmerje verjetnosti, da dogodek zgodi in da se ne zgodi:

$$\frac{p(x)}{1 - p(x)}$$

■ Zavzame vrednosti med 0 in $+$ neskončnostjo

Preprosta logistična regresija

- Logaritmiranje \Rightarrow med $-$ in $+$ neskončnostjo

$$\ln \left[\frac{p(x)}{1-p(x)} \right] = \log it$$

Populacija:

$$\ln \left[\frac{\pi(x)}{1-\pi(x)} \right] = a + bx = \log it$$

Ocena parametrov iz vzorca:

$$\ln \left[\frac{p(x)}{1-p(x)} \right] = a + bx = \log it$$

$$p(x) = \frac{e^{a+bx}}{1+e^{a+bx}}$$

Verjetnost, da dogodek
zgodil pri danem X

Pomen regresijskega koef. b

□ Lin. regresija: $b = y_{(x+1)} - y_{(x)}$ $\ln\left[\frac{p(x)}{1-p(x)}\right] = a + bx = \log it$

□ Log. regresija: $b = \log it_{(x+1)} - \log it_{(x)}$

□ Binarni X (0/1) $\Rightarrow x=0, x+1=1.$

□ Obet da bo $y=1$:

$$O_{x=1} = \frac{p(1)}{1-p(1)}$$

$$O_{x=0} = \frac{p(0)}{1-p(0)}$$

□ Logaritmiranje:

$$\log it_{x=1} = \ln\left[\frac{p(1)}{1-p(1)}\right]$$

$$\log it_{x=0} = \ln\left[\frac{p(0)}{1-p(0)}\right]$$

Pomen regresijskega koef. b

- Log. razmerja obetov:

$$RO = \frac{O_{x=1}}{O_{x=0}} = \frac{\left[\frac{p(1)}{1-p(1)} \right]}{\left[\frac{p(0)}{1-p(0)} \right]}$$

$\leftarrow O_{x=1} = \frac{p(1)}{1-p(1)}$
 $\leftarrow O_{x=0} = \frac{p(0)}{1-p(0)}$

$$\ln RO = \ln \frac{\left[\frac{p(1)}{1-p(1)} \right]}{\left[\frac{p(0)}{1-p(0)} \right]} = \log it(1) - \log it(0)$$

$$\ln RO = b = \log it_{(1)} - \log it_{(0)}$$

$$RO = e^b$$

Ocenjevanje b v vzorcu

- $p(x)? \Rightarrow b$ in a ? $\ln \left[\frac{p(x)}{1-p(x)} \right] = a + bx = \log it$
- Lin. regresija: metoda najmanjših kvadratov ostankov.
- Log. regresija: metoda največjega verjetja (maximum likelihood method).
 - Funkcija največjega verjetja
 - Oz. logaritem funkcije verjetja ("log likelihood"):
 - Nelinearna funkcija parametrov modela a in b .
 - Iteracijska metoda, več ocen parametrov. Nove ocene, dokler še zveča funkcijo največjega verjetja. Lokalna/globalna točka največjega verjetja. Start?

Ocenjevanje b v populaciji

- Vzorčna porazdelitev b-ja: normalna porazdelitev = > interval zaupanja:

$$\beta = b \pm z * SE_{(b)}$$

- Vzorčna ocena razmerja obetov: nenormalna, nesimetrična porazdelitev
- Interval zaupanja za RO iz b (spodnjo/zgornjo mejo):

$$e^{b \pm z * SE_b}$$

Multipla logistična regresija

- Matematični model:

$$p(x) = \frac{e^{a+b_1x_1+b_2x_2+\dots+b_kx_k}}{1+e^{a+b_1x_1+b_2x_2+\dots+b_kx_k}}$$

- $b(k)$: sprememba v logitu, ki spremlja spremembo $X(k)$ za 1 enoto, medtem ko se ostale sprem. ne spreminjajo.
 - Regresijski koef. b = logaritem razmerja obolevnosti
 - Razmerje obolevnosti = antilogaritem b : e^b
 - Metoda največjega verjetja
-

Testiranje pomembnosti modela kot celote

- Ali se log verjetja modela s spremenljivkami X statistično značilno poveča v primerjavi z log verjetja modela brez njih.
- G- statistika- test razmerja dveh verjetij:

$$G = 2 \ln \left[\frac{\text{verjetje}_{SPREM.X}}{\text{verjetje}_0} \right]$$

- Porazdeljuje po hi- kvadrat, df: k-1
 - k: št. vseh spremenljivk v modelu
 - V SPSS-u: "Omnibus test"
-

Vrednotenje prileganja modela

- Nagelkerkejev R^2

 - Test hi- kvadrat

 - Hosmer- Lemeshov test:
 - Oba primerjata opazovano število enot, pri katerem se je opazovani dogodek zgodil, s pričakovanim številom, ki temelji na enačbi logistične regresije.
-

Testiranje pomembnosti b

□ $H_0: b=0, H_a: b \neq 0$

□ Testi:

■ Z:
$$z = \frac{b-0}{SE_b}$$

■ Waldova statistika:

$$\chi^2 = \left(\frac{b}{SE_b} \right)^2$$

■ G- statistika:

$$G = 2 \ln \left[\frac{\text{verjetje}_x}{\text{verjetje}_0} \right]$$

G- statistika:
Model, ki
vključuje X, pove
več o Y.

Napovedni dejavniki nenamenske uporabe zdravil

Nagelkerke R. Square = 0.185

Hosmer Lemeshow Test: Chi square = 5.240, df = 7, p = 0.631

Method: "Enter"

Variables in the model	B	S.E.	Wald	df	Sig.	Exp(B)	95% C.I. for EXP(B)	
							Lower Bound	Upper Bound
(Constant)	-0.933	0.541	2.973	1	0.085	0.394		
Specialization (no => yes)	1.183	0.531	4.961	1	0.026	3.263	1.152	9.241
Institution (public => private)	-1.786	0.801	4.979	1	0.026	0.168	0.035	0.805
Region (Ljubljana & Maribor => Other regions)	-0.942	0.472	3.975	1	0.046	0.390	0.154	0.984
Review frequency of treatment novelties (monthly => weekly)	-0.763	0.478	2.550	1	0.110	0.466	0.183	1.190

Characteristics that predict physicians' awareness of their off-label prescribing. Logistic regression model with statistically significant predictors.

Opisne sprem.v regresiji

- Opisne sprem. dve vrednosti => kodiranje 0/1
- Linearna regresija:
 - Enačbi za pričakovani povprečni vrednosti Y:

$$\bar{y}_{x=0} = a + b_1 * 0 = a$$

$$\bar{y}_{x=1} = a + b_1 * 1 = a + b_1$$

- Podobno analogija tudi za log. regresijo
-

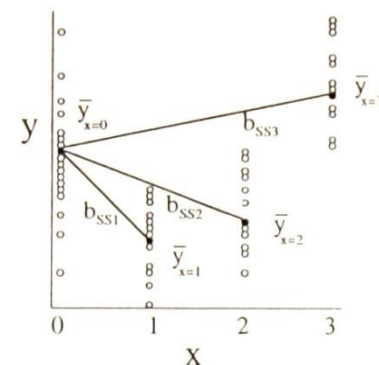
Slepe spremenljivke ("dummy var.")

- Opisne sprem. >2 vrednosti
- Podobno kot pri binarni sprem.?!
- Nove sprem.: vrednosti 0/1

Preprosto kodiranje

- Primerjalna kategorija (nima svoje SS)
- $b(ss)$: razlika med aritm. sredino \bar{Y} pri opazovani kategoriji in aritm. sredino primerjalna kategorije oz. konstanto a
- Za nominalne opisne sprem. z več kategorijami

	slepe spremenljivke			enačba
kategorija osn. sprem.	SS1	SS2	SS3	
0	0	0	0	$\bar{y}_0 = a$
1	1	0	0	$\bar{y}_1 = a + b_{SS1}$
2	0	1	0	$\bar{y}_2 = a + b_{SS2}$
3	0	0	1	$\bar{y}_3 = a + b_{SS3}$
reg. koef.	b_{SS1}	b_{SS2}	b_{SS3}	



Slika 12

Sekvenčno kodiranje:

- Osnovna kategorija (nima svoje SS)
 - Ni več vedno primerjalna
- $b(ss)$: razlika med aritm. sredino Y pri opazovani kategoriji in aritm. sredino prejšnje kategorije
- Za opisne sprem. ordinalnega tipa

kategorija	slepe spremenljivke			enačba
	SS1	SS2	SS3	
osn. sprem.				
0	0	0	0	$\bar{y}_0 = a$
1	1	0	0	$\bar{y}_1 = a + b_{SS1}$
2	1	1	0	$\bar{y}_2 = a + b_{SS1} + b_{SS2}$
3	1	1	1	$\bar{y}_3 = a + b_{SS1} + b_{SS2} + b_{SS3}$
reg. koef.	b_{SS1}	b_{SS2}	b_{SS3}	

Slika 13

